

Project: ZuluNation Motor

Introduction /Overview:

The HR department at Zulunation Motors wants to take some initiatives to improve employee satisfaction levels at the company. They collected data from employees, and They have the following question:

what's likely to make the employee leave the company?

Data Source:

The dataset is a fictitious example created for practice and knowledge.

Objective:

The goals in this project are to analyze the data collected by the HR department and to build a model/s that predicts whether or not an employee will leave the company. By successfully predicting which employees are likely to quit, it might be possible to identify factors that contribute to their decision to leave.

Setting up the environment, importing packages and load the dataset :

```
In [84]: 1 import pandas as pd  
2 import numpy as np
```

```
In [2]: 1 data = pd.read_csv('C:/Users/engmo/OneDrive/Desktop/Projects to show of
```

EDA and Data cleaning

```
In [3]: 1 data.head()
```

```
Out[3]:
```

	satisfaction_level	last_evaluation	number_project	average_monthly_hours	time_spend_company
0	0.38	0.53	2	157	
1	0.80	0.86	5	262	
2	0.11	0.88	7	272	
3	0.72	0.87	5	223	
4	0.37	0.52	2	159	

In [4]: 1 data.columns

Out[4]: Index(['satisfaction_level', 'last_evaluation', 'number_project',
'average_monthly_hours', 'time_spend_company', 'Work_accident', 'left',
'promotion_last_5years', 'Department', 'salary'],
dtype='object')

In [21]: 1 data.rename(columns={'average_monthly_hours':'average_monthly_hours'},in

In [5]: 1 data.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 14999 entries, 0 to 14998
Data columns (total 10 columns):
#   Column                Non-Null Count  Dtype
---  -
0   satisfaction_level     14999 non-null  float64
1   last_evaluation        14999 non-null  float64
2   number_project         14999 non-null  int64
3   average_monthly_hours  14999 non-null  int64
4   time_spend_company     14999 non-null  int64
5   Work_accident          14999 non-null  int64
6   left                   14999 non-null  int64
7   promotion_last_5years  14999 non-null  int64
8   Department             14999 non-null  object
9   salary                 14999 non-null  object
dtypes: float64(2), int64(6), object(2)
memory usage: 1.1+ MB
```

In [6]:

```
1  ## Lets standarize the column names
2  data.columns = [col.lower()for col in data.columns]
3  ## Lets rename some of the columns to improve simpilicity.
4  data.rename(columns={'time_spend_company':'tenure'},inplace=True)
5  ## display the first 10 rows of the dataset
6  data.head(10)
```

Out[6]:

	satisfaction_level	last_evaluation	number_project	average_monthly_hours	tenure	work_acc
0	0.38	0.53	2	157	3	
1	0.80	0.86	5	262	6	
2	0.11	0.88	7	272	4	
3	0.72	0.87	5	223	5	
4	0.37	0.52	2	159	3	
5	0.41	0.50	2	153	3	
6	0.10	0.77	6	247	4	
7	0.92	0.85	5	259	5	
8	0.89	1.00	5	224	5	
9	0.42	0.53	2	142	3	

```
In [7]: 1 # A descriptive analysis of the data.
        2 data.describe()
```

```
Out[7]:
```

	satisfaction_level	last_evaluation	number_project	average_monthly_hours	tenure
count	14999.000000	14999.000000	14999.000000	14999.000000	14999.000000
mean	0.612834	0.716102	3.803054	201.050337	3.498233
std	0.248631	0.171169	1.232592	49.943099	1.460136
min	0.090000	0.360000	2.000000	96.000000	2.000000
25%	0.440000	0.560000	3.000000	156.000000	3.000000
50%	0.640000	0.720000	4.000000	200.000000	3.000000
75%	0.820000	0.870000	5.000000	245.000000	4.000000
max	1.000000	1.000000	7.000000	310.000000	10.000000

```
In [8]: 1 data['promotion_last_5years'].value_counts(normalize=True)
```

```
Out[8]: 0    0.978732
        1    0.021268
        Name: promotion_last_5years, dtype: float64
```

Check missing values:

```
In [9]: 1 data.isna().sum() #TO find the missing values in the dataset
```

```
Out[9]: satisfaction_level      0
        last_evaluation         0
        number_project          0
        average_monthly_hours   0
        tenure                  0
        work_accident           0
        left                    0
        promotion_last_5years    0
        department              0
        salary                  0
        dtype: int64
```

Check duplicates:

```
In [20]: 1 duplicates = data.duplicated() ## Inspect duplicates.
2         duplicates_rows = data[duplicates] ## Creating a DataFrame For duplicates
3         duplicates_rows.head() ## showing the first 5 rows of the duplicates rows
```

```
Out[20]:
```

	satisfaction_level	last_evaluation	number_project	average_monthly_hours	tenure	work_level
396	0.46	0.57	2	139	3	
866	0.41	0.46	2	128	3	
1317	0.37	0.51	2	127	3	
1368	0.41	0.52	2	132	3	
1461	0.42	0.53	2	142	3	

Check outliers:

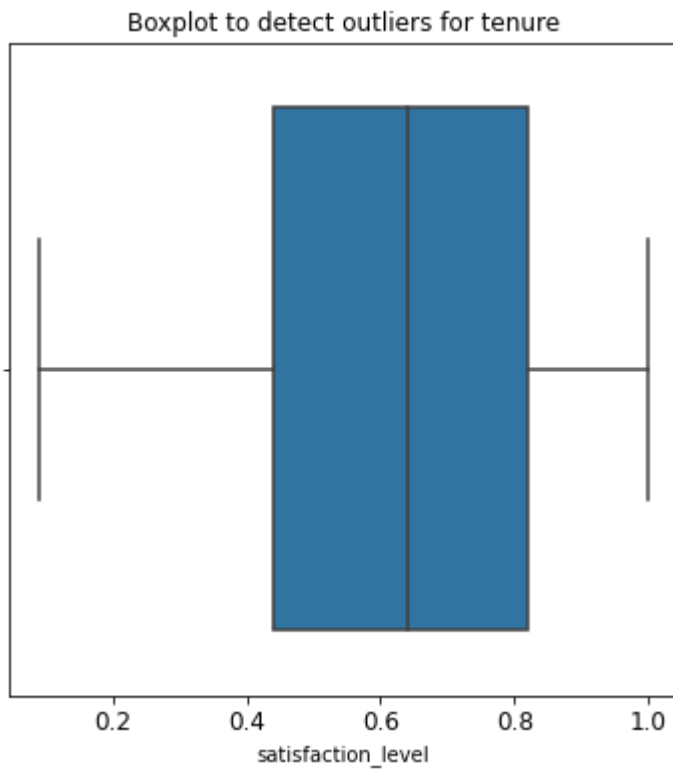
```
In [29]: 1 ## Let's check outliers in our dataset:
2         ## Lets bring the descriptive analysis of our data to start
3         data.describe()
```

```
Out[29]:
```

	satisfaction_level	last_evaluation	number_project	average_monthly_hours	tenure
count	14999.000000	14999.000000	14999.000000	14999.000000	14999.000000
mean	0.612834	0.716102	3.803054	201.050337	3.498230
std	0.248631	0.171169	1.232592	49.943099	1.460130
min	0.090000	0.360000	2.000000	96.000000	2.000000
25%	0.440000	0.560000	3.000000	156.000000	3.000000
50%	0.640000	0.720000	4.000000	200.000000	3.000000
75%	0.820000	0.870000	5.000000	245.000000	4.000000
max	1.000000	1.000000	7.000000	310.000000	10.000000

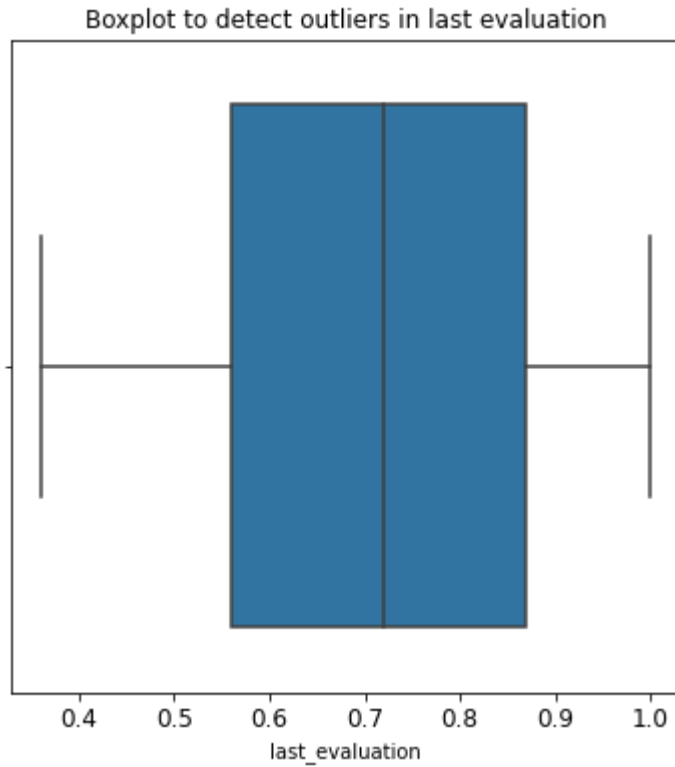
```
In [32]: 1  ## Let's create a boxplot to visulaize the outliers
2  ## Satisfaction Level
3  import matplotlib.pyplot as plt
4  import seaborn as sns
5  plt.figure(figsize=(6,6))
6  plt.title('Boxplot to detect outliers in satisfaction level', fontsize
7  plt.xticks(fontsize=12)
8  plt.yticks(fontsize=12)
9  sns.boxplot(x=data['satisfaction_level'])
10 plt.show
```

Out[32]: <function matplotlib.pyplot.show(close=None, block=None)>



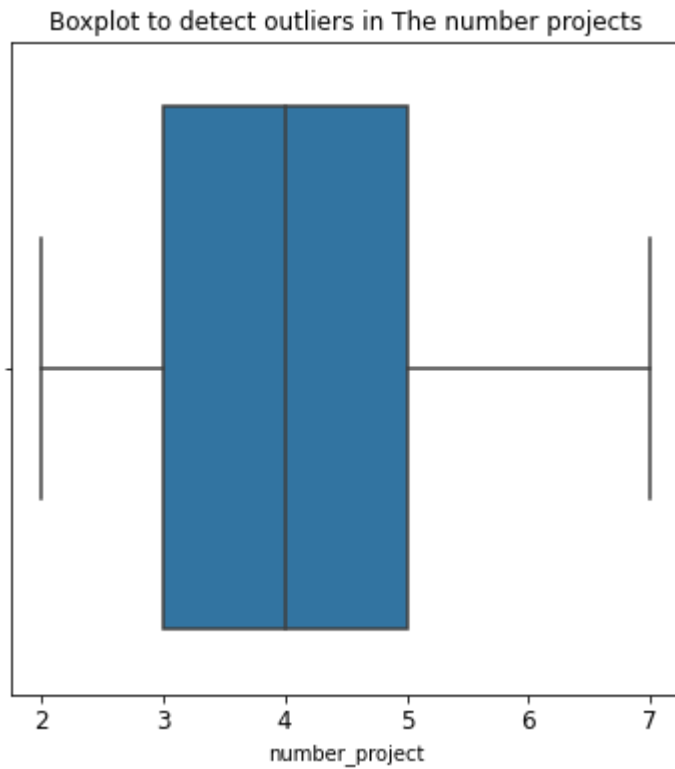
```
In [34]: 1  ## Lets visulaize the last evaluation
          2
          3  plt.figure(figsize=(6,6))
          4  plt.title('Boxplot to detect outliers in last evaluation', fontsize=12)
          5  plt.xticks(fontsize=12)
          6  plt.yticks(fontsize=12)
          7  sns.boxplot(x=data['last_evaluation'])
          8  plt.show
```

Out[34]: <function matplotlib.pyplot.show(close=None, block=None)>



```
In [35]: 1  ## Lets visulaize number_project
2  plt.figure(figsize=(6,6))
3  plt.title('Boxplot to detect outliers in The number projects',fontsize=
4  plt.xticks(fontsize=12)
5  plt.yticks(fontsize=12)
6  sns.boxplot(x=data['number_project'])
7  plt.show
```

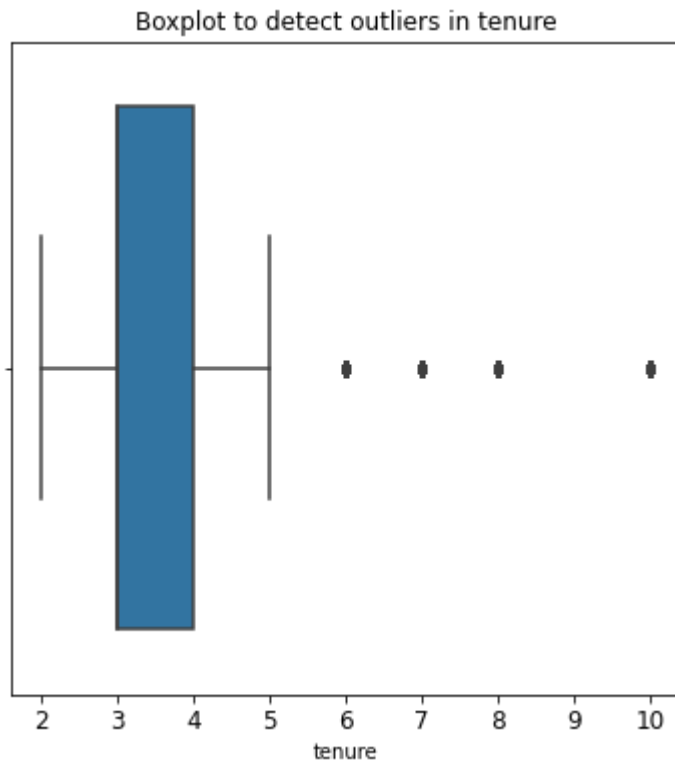
Out[35]: <function matplotlib.pyplot.show(close=None, block=None)>



```
In [ ]: 1  ## Lets visulaize average monthly hours
2  plt.figure(figsize=(6,6))
3  plt.title('Boxplot to detect outliers in The average monthly hours',font
4  plt.xticks(fontsize=12)
5  plt.yticks(fontsize=12)
6  sns.boxplot(x=data['average_monthly_hours'])
7  plt.show
```

```
In [37]: 1  ## Lets visulaize tenure
2  plt.figure(figsize=(6,6))
3  plt.title('Boxplot to detect outliers in tenure',fontsize=12)
4  plt.xticks(fontsize=12)
5  plt.yticks(fontsize=12)
6  sns.boxplot(x=data['tenure'])
7  plt.show
```

Out[37]: <function matplotlib.pyplot.show(close=None, block=None)>



The boxplot above shows that there are outliers oin the `tenure` variable.

It would be helpful to investigate how many rows in the data contain outliers in the `tenure` column.

In [40]:

```

1  ##Let's determine the number of rows containing outliers.
2  ## Let's start by computing IQR(interquartile range)
3
4  percentile25 = data['tenure'].quantile(0.25)
5  percentile75 = data['tenure'].quantile(0.75)
6
7  iqr = percentile75 - percentile25
8
9  ## Lets define the upper limit and the lower limit for non-outliers val
10 upper_limit = percentile75 + 1.5 * iqr
11 lower_limit = percentile25 - 1.5* iqr
12 print("lower Limit:", lower_limit)
13 print("Upper Limit:", upper_limit)
14
15 ## Let's identify the subset of the data containinf outliers in `tenure
16 outliers = data[(data['tenure'] > upper_limit) | (data['tenure'] < lowe
17
18 ## Let's identify how many rows containing outliers
19 print("Number of rows in the data containing outliers in tenure:",len(o
20

```

lower Limit: 1.5

Upper Limit: 5.5

Number of rows in the data containing outliers in tenure: 1282

Let's keep the outliers for now until determining which model to use and understand the model sensitivity to these outliers

Continuing EDA :

In [43]:

```

1  ## Let's undersrtand how many employees left and what their representat
2  print(data['left'].value_counts())
3  print(data['left'].value_counts(normalize=True)*100)

```

0 11428

1 3571

Name: left, dtype: int64

0 76.191746

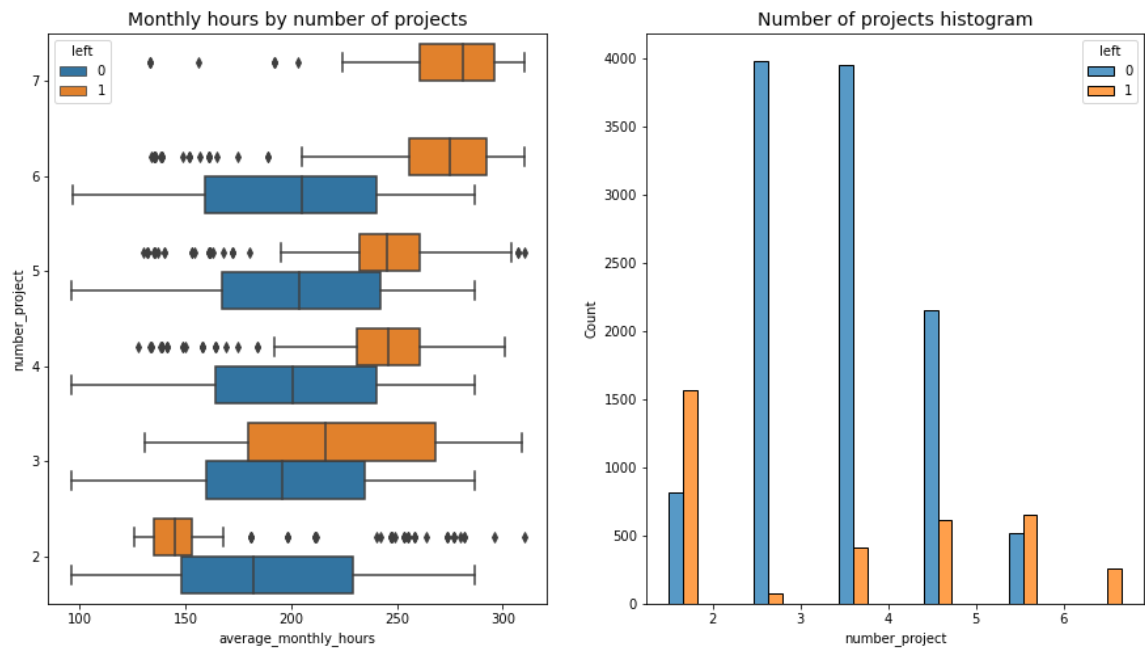
1 23.808254

Name: left, dtype: float64

Data visualizations:

```
In [61]: 1  ## Let's create a boxplot showing 'average_monthly_hours' distributions
2  fig, ax = plt.subplots(1,2,figsize=(15,8))
3  sns.boxplot(data=data,x='average_monthly_hours',y='number_project',hue=
4  ax[0].invert_yaxis()
5  ax[0].set_title('Monthly hours by number of projects', fontsize=14)
6  ## Now let's creat a histogram to visualize the distribution of number_
7  sns.histplot(data=data,x='number_project',hue='left', multiple='dodge',
8  ax[1].set_title("Number of projects histogram", fontsize =14)
9
10 plt.show
11
```

```
Out[61]: <function matplotlib.pyplot.show(close=None, block=None)>
```



Notes:

It's natural that people work on more projects tend to work longer hours, this appears to be the case here, however a few things stand out from this plot.

1. There are two groups of employees who left the company: (A) those who worked considerably less than their peers with the same number of projects, and (B) who worked much more. group (A) might be the people that who are serving their contract notice period and they were assigned to fewer hours.
2. The optimal number of projects for employees to work on seems to be 3-4, the ratio of left/stayed is very small.
3. If the employee should work 40 hours/week and 166.67 hours/month, the mean average of monthly hours is 201 and some of the employees worked 301 hours, it seems that the employees are overworked.
4. Every employee with project more than 6 left the company.

```
In [62]: 1 ## Let confirm if all the employees with 7 projects left the company.
2
3 data[data['number_project'] == 7]['left'].value_counts()
```

```
Out[62]: 1    256
Name: left, dtype: int64
```

This confirms that all employees with 7 projects left the company

```
In [65]: 1 ##Let's examine the average monthly hours Vs the employee satisfaction
2
3 #Let's create scatterplot of the avergae monthly hours vs the employee
4 plt.figure(figsize=(12,6))
5 sns.scatterplot(data=data,x='average_monthly_hours',y='satisfaction_level')
6 plt.axvline(x=166.67,color = '#ff6361',label='166.67 hrs./mo.',ls='--')
7 plt.legend(labels=['166.67 hrs./mo.','left','stayed'])
8 plt.title('Monthly hours by satisfaction score',fontsize=14)
9 plt.show
```

```
Out[65]: <function matplotlib.pyplot.show(close=None, block=None)>
```



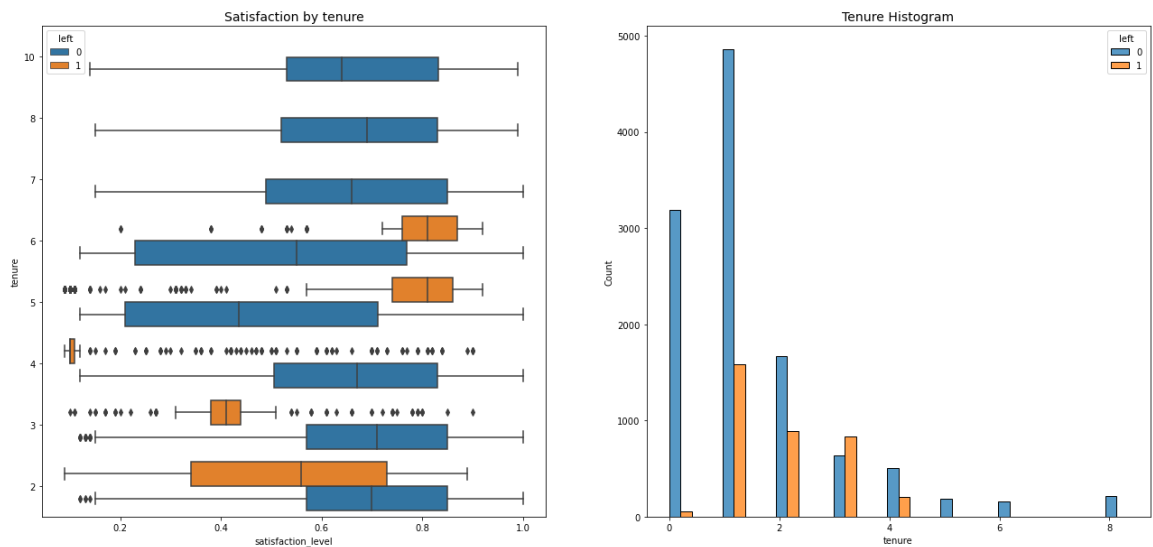
Notes:

1. From the above plot notice that a group of employees worked between 230 -330 hours/month and this is more that the average working hours, this could be the reason for their satisfaction level.
2. The plot also shows that there is a gorup of employees who worked minimum hours comparing to their peers and yet they left, their satisfaction score is around 0.4.
3. Finally, a third group who have worked between 210 - 280 hours/monnth they have left but their satisfaction level is above 0.75.

```

In [71]: 1  ## Let's visualize the satisfaction level vs the tenure.
2
3  #Lets set the figure axis
4  fig, ax = plt.subplots(1,2, figsize=(22,10))
5
6  ## Lets create a boxplot showing the distributions of the satisfaction
7
8  sns.boxplot(data=data,x="satisfaction_level",y='tenure',hue='left',orie
9  ax[0].invert_yaxis()
10 ax[0].set_title('Satisfaction by tenure',fontsize=14)
11
12
13 ##Lets create a histogram showing the distribution of tenure, comparing
14 sns.histplot(data=data,x='tenure',hue='left',multiple='dodge',shrink=5,
15 ax[1].set_title('Tenure Histogram',fontsize='14')
16
17 plt.show()
18

```



Observations:

1. Employees with longer-tenure tends to stay and the have the same satisfaction level as those who newly joined the company.
2. Employees at 4 years tenure have unusual satisfaction score, it worth checking the company policies or any changes happened at 4 year mark.
3. The majority of employees who left worked few years and they have low satisfaction level.

```

In [72]: 1  ## Let's calculate the mean a median satisfaction scores of employees w
2  data.groupby(['left'])['satisfaction_level'].agg([np.mean,np.median])
3

```

```

Out[72]:      mean  median
left
0  0.666810  0.69
1  0.440098  0.41

```

Observations: The mean and median for those who left are lower than the score of the employees who stayed, Among the employees who stayed the mean is lower than the median which indicates that the satisfaction scores among those who stayed are skewed to the left.

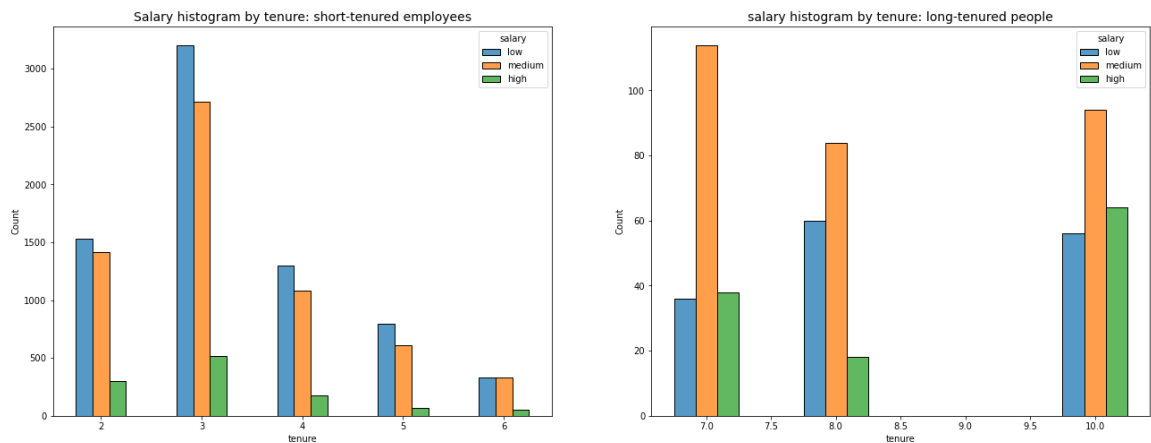
In [78]:

```

1  ## Let's examine the salary levels for different tenures.
2  # let's set the figure and axes
3  fig, ax = plt.subplots(1,2,figsize=(22,8))
4  # Lets define short-tenured employees.
5  tenure_short = data[data['tenure'] < 7]
6  tenure_long = data[data['tenure'] > 6]
7
8  #Let's plot short_tenured histogram:
9  sns.histplot(data=tenure_short,x = 'tenure',hue='salary', hue_order=['1
10 ax[0].set_title('Salary histogram by tenure: short-tenured employees',f
11
12 #And the Long_tenured employees:
13 sns.histplot(data=tenure_long,x='tenure', shrink=.5,hue='salary', hue_o
14 discrete=1,ax=ax[1])
15 ax[1].set_title('salary histogram by tenure: long-tenured people',font
16 plt.show

```

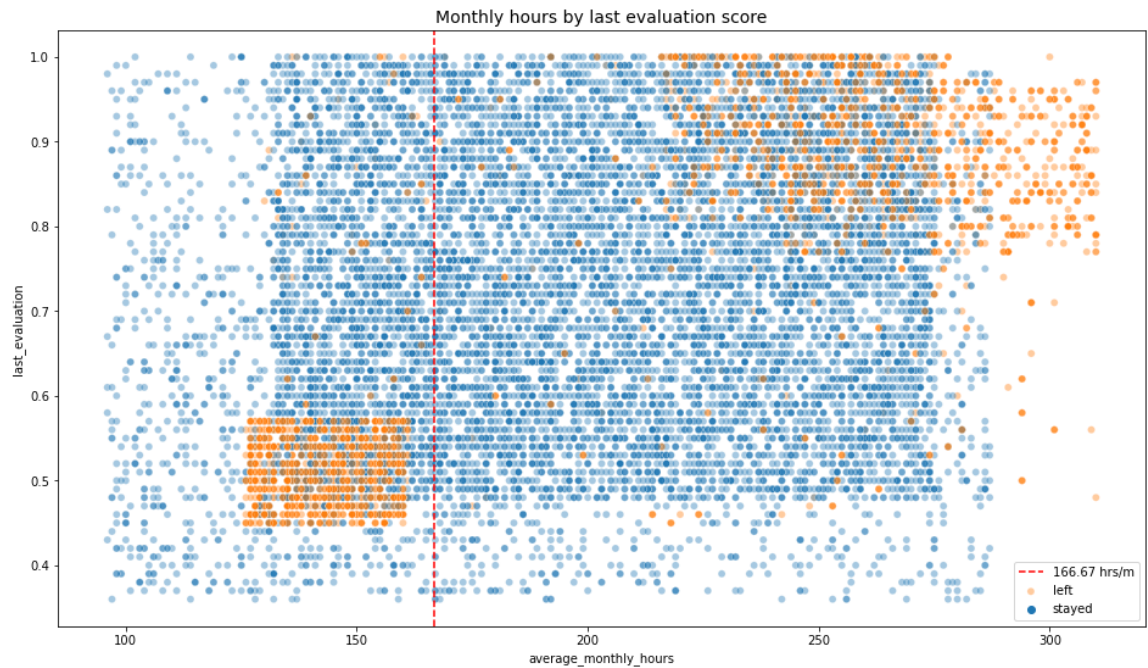
Out[78]: <function matplotlib.pyplot.show(close=None, block=None)>



Observation: Being a long-tenured employee does not necessarily correlate with having a higher salary.

```
In [79]: 1  ## Let's Look at the average monthly hours vs the evaluation scores:
2
3  plt.figure(figsize=(16,9))
4  sns.scatterplot(data=data,x='average_monthly_hours',y='last_evaluation')
5  plt.axvline(x=166.67,color='red',label='166.67 hr/m',ls='--')
6  plt.legend(labels=['166.67 hrs/m', 'left', 'stayed'])
7  plt.title('Monthly hours by last evaluation score',fontsize=14)
8  plt.show
```

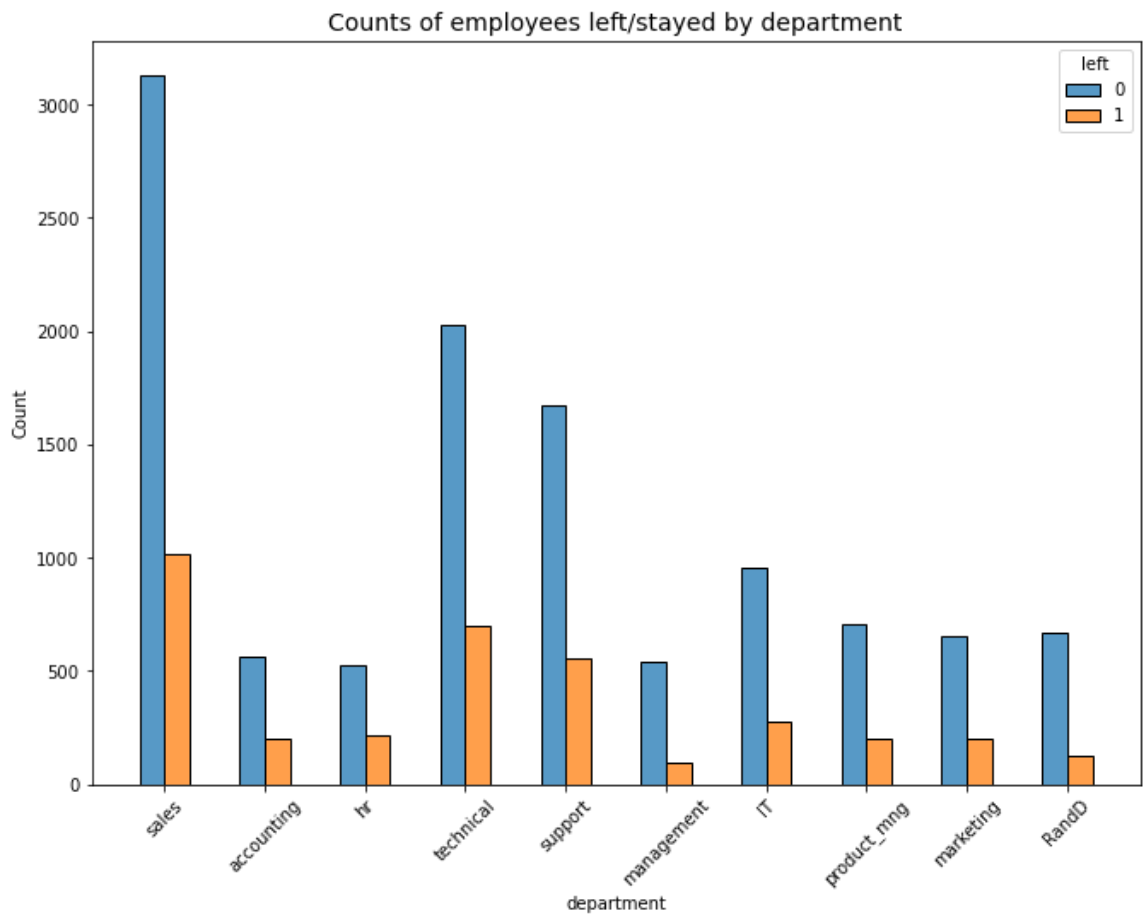
```
Out[79]: <function matplotlib.pyplot.show(close=None, block=None)>
```



Observations: 1. There seems to be a correlation between hours worked and evaluation score. 2. Most of the employees in this company work well over 167 hours per month.

```
In [80]: 1 ## Now Lets visualize the employees who left versus their department:  
2 plt.figure(figsize=(11,8))  
3 sns.histplot(data=data,x='department',hue='left',discrete=1,hue_order=[  
4  
5 plt.xticks(rotation=45)  
6 plt.title('Counts of employees left/stayed by department',fontsize=14)
```

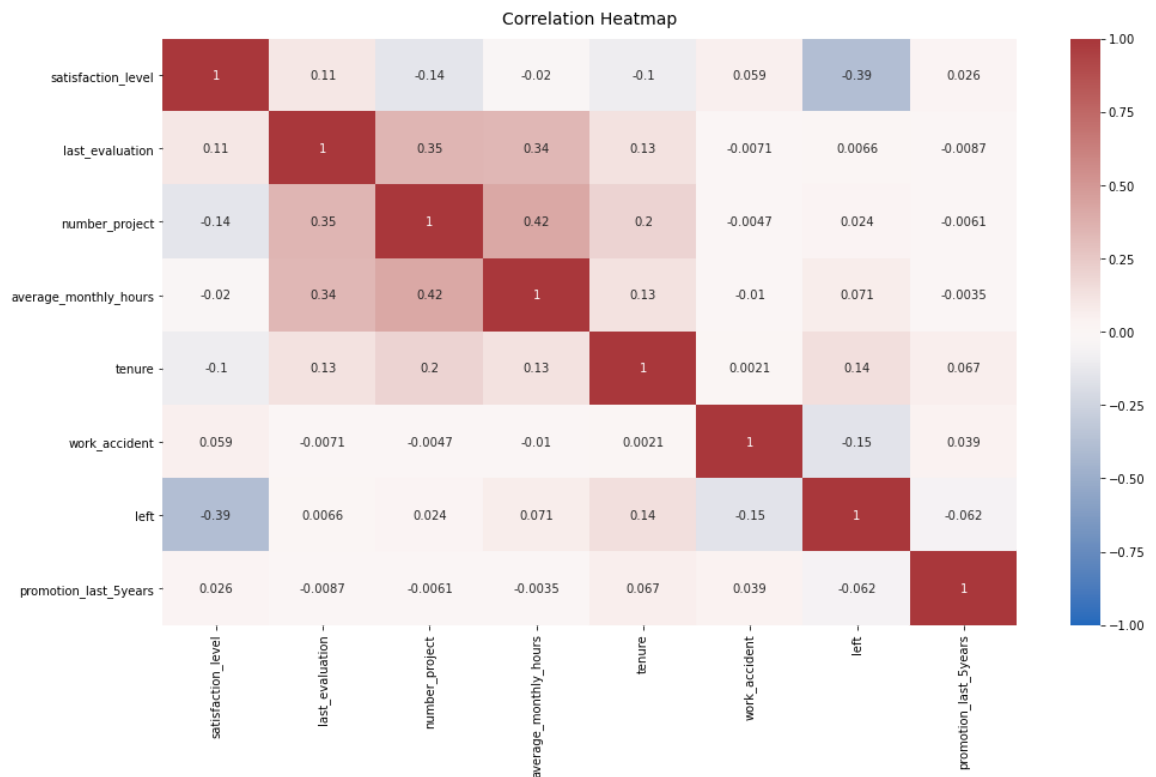
Out[80]: Text(0.5, 1.0, 'Counts of employees left/stayed by department')



Observation: There doesn't seem to be any department that differs significantly in its proportion of employees who left to those who stayed

```
In [83]: 1 ## Let's create a heatmap to check for a strong correlation between var
2 plt.figure(figsize=(16,9))
3 heatmap= sns.heatmap(data.corr(),vmin=-1,vmax=1,annot=True,cmap=sns.col
4 heatmap.set_title("Correlation Heatmap",fontdict={'fontsize':14},pad=12
5 plt.show
```

```
Out[83]: <function matplotlib.pyplot.show(close=None, block=None)>
```



Observation: The correlation heatmap confirms that the number of projects, monthly hours, and evaluation scores all have some positive correlation with each other, and whether an employee leaves is negatively correlated with their satisfaction level.

Insights:

1. Leaving is tied to longer working hours, many projects, and generally lower satisfaction levels.
2. There's a sizeable group of employees at this company who are probably burned out.
3. It can be ungratifying to work long hours and not receive promotions or good evaluation scores.
4. It also appears that if the employee cross the 6 years tenure mark they tend to stay.

Notes: By examining the EDA insights and outcomes, we can start by choosing and developing the model.

Model Development:

Since the outcome variable is categorical, Lets develop a logisitc regression model and dicison tree model as well and compare how they performed.

Before splitting the data, lets encode the nonnumerical variables in the dataset, department and salary

Approach (A): Logistic Regression

```
In [88]: 1  ##Lets copy the dataframe.
          2  data_new = data.copy()
          3
          4  ##Notice that salary is categorical but its not ordinal, there is a hie
          5
          6  data_new['salary'] = (data_new['salary']).astype('category').cat.set_cat
          7
          8  ## And Lets dummy the department for modeling.
          9  data_new= pd.get_dummies(data_new,drop_first=False)
         10  data_new.head()
```

```
Out[88]:
```

	satisfaction_level	last_evaluation	number_project	average_monthly_hours	tenure	work_ac
0	0.38	0.53	2	157	3	
1	0.80	0.86	5	262	6	
2	0.11	0.88	7	272	4	
3	0.72	0.87	5	223	5	
4	0.37	0.52	2	159	3	

```
In [89]: 1  ## Since logistic regression is sensitive to outliers, Lets remove the
          2  data_new = data_new[(data_new['tenure'] >= lower_limit) & (data_new['te
          3
          4  data_new.head()
```

```
Out[89]:
```

	satisfaction_level	last_evaluation	number_project	average_monthly_hours	tenure	work_ac
0	0.38	0.53	2	157	3	
2	0.11	0.88	7	272	4	
3	0.72	0.87	5	223	5	
4	0.37	0.52	2	159	3	
5	0.41	0.50	2	153	3	

```
In [90]: 1  ## Now Lets Isolate the outcome variable and assign it to y.
          2
          3  y=data_new['left']
          4  y.head()
```

```
Out[90]: 0    1
          2    1
          3    1
          4    1
          5    1
          Name: left, dtype: int64
```

```
In [91]: 1 ##Now Lets select the features and assign it to X.
2
3 X = data_new.drop('left',axis=1)
4 X.head()
```

```
Out[91]:
```

	satisfaction_level	last_evaluation	number_project	average_monthly_hours	tenure	work_ac
0	0.38	0.53	2	157	3	
2	0.11	0.88	7	272	4	
3	0.72	0.87	5	223	5	
4	0.37	0.52	2	159	3	
5	0.41	0.50	2	153	3	

```
In [98]: 1 ## Here we would split hte data into training set and testing test, we
2 ## testing size would be 25%.
3 ## Lets import the required packages.
4 from sklearn.model_selection import train_test_split
5
6
7 X_train,X_test,y_train,y_test = train_test_split(X,y,test_size=0.25,str
8
```

```
In [99]: 1 ## Let's Construct the Logistice regression.
2 ## Lets Import the required packages.
3 from sklearn.linear_model import LogisticRegression
4 log = LogisticRegression(random_state=42,max_iter = 500)
5
6 ## and fit the regression model
7
8 Log_f = log.fit(X_train,y_train)
```

```
In [101]: 1 ## Lets test the model by using the model to get predictions.
2 y_pred = Log_f.predict(X_test)
```

```
In [ ]: 1 ## Lets create a confusion matrix to visualize the results.
2 ## Let's import the required packages (PS. Lets import all the packages
3
4 from sklearn.metrics import accuracy_score, precision_score, recall_sco
5 f1_score, confusion_matrix, ConfusionMatrixDisplay, classification_repo
6 from sklearn.metrics import roc_auc_score, roc_curve
7
8 cm = confusion_matrix(y_test,y_pred,labels=Log_f.classes_)
9 cm_disp = ConfusionMatrixDisplay(confusion_matrix=cm,display_labels=Log
10 cm_disp.plot(values_format='')
11
12 plt.show()
```

Notice: The number of false negatives is higher than the number of false positives. This implies that the model is more conservative in predicting the positive class; it's more likely to miss positives (predict them as negatives, employees stayed), Lets check the outcome variable balanced in dataset

```
In [104]: 1 data_new['left'].value_counts(normalize=True)*100
```

```
Out[104]: 0    75.490268
          1    24.509732
          Name: left, dtype: float64
```

There is an approx. 75% to 25% split, which its not severely imbalance. lets continue evaluating the model.

```
In [105]: 1 ## Lets create a classification report
          2 target_names = ['Predicated would not leave', 'Predicated would leave']
          3 print(classification_report(y_test,y_pred,target_names=target_names))
```

	precision	recall	f1-score	support
Predicated would not leave	0.86	0.91	0.88	2589
Predicated would leave	0.66	0.56	0.60	841
accuracy			0.82	3430
macro avg	0.76	0.73	0.74	3430
weighted avg	0.81	0.82	0.81	3430

Observations:

1. The model is quite good at predicting the employees who would not leave (Class 0), as indicated by high precision, recall, and F1-score for this class.
2. The model struggles relatively more with predicting the employees who would leave (Class 1), which is evident from the lower recall and F1-score.
3. Improving the model could involve addressing the imbalance, perhaps by resampling the dataset

Approach (B): Tree-based Model

```
In [107]: 1  ## Lets Isolate, build, fit and evaluate a decision tree model.
2  ## Prepare the dataset
3  ##Lets copy the dataframe.
4  data_two = data.copy()
5
6  ##Notice that salary is categorical but its not ordinal, there is a hie
7
8  data_two['salary'] = (data_two['salary']).astype('category').cat.set_cat
9
10 ## And Lets dummy the department for modeling.
11 data_two= pd.get_dummies(data_two,drop_first=False)
12 data_two.head()
```

```
Out[107]:
```

	satisfaction_level	last_evaluation	number_project	average_monthly_hours	tenure	work_ac
0	0.38	0.53	2	157	3	
1	0.80	0.86	5	262	6	
2	0.11	0.88	7	272	4	
3	0.72	0.87	5	223	5	
4	0.37	0.52	2	159	3	

```
In [108]: 1  ## isolate, build and fit.
2  y = data_two['left']
3  y.shape
```

```
Out[108]: (14999,)
```

```
In [111]: 1  #Features.
2  X= data_two.drop('left',axis=1)
3  X.head()
```

```
Out[111]:
```

	satisfaction_level	last_evaluation	number_project	average_monthly_hours	tenure	work_ac
0	0.38	0.53	2	157	3	
1	0.80	0.86	5	262	6	
2	0.11	0.88	7	272	4	
3	0.72	0.87	5	223	5	
4	0.37	0.52	2	159	3	

```
In [112]: 1  ## here we would split the data into training, validating and testing s
2  X_train,X_test,y_train,y_test = train_test_split(X,y,test_size=0.25,ran
3
```

```
In [113]: 1 ##Lets import the necessary packages to build a tree based models.
2 from sklearn.tree import DecisionTreeClassifier
3 from sklearn.ensemble import RandomForestClassifier
4 from sklearn.model_selection import GridSearchCV, train_test_split
5 from sklearn.tree import plot_tree
```

```
In [117]: 1 ## Instantiate model & setting up cross-validated grid-search to search
2 tree= DecisionTreeClassifier(random_state=0)
3 ##Assigning a dictionary of hyperparameters to search over
4
5 cv_params = {'max_depth':[4,6,8,None], 'min_samples_leaf':[2,5,1], 'min_s
6
7 ##Assigning the scoring metrics
8
9 scoring = {'accuracy', 'precision', 'recall', 'f1', 'roc_auc'}
10
11 ## Insintiate the GridSearch
12 tree1 = GridSearchCV(tree,cv_params,scoring=scoring,cv=5,refit='roc_auc
13
```

Fit the Tree on the training data.

```
In [118]: 1 %%time
2 tree1.fit(X_train,y_train)
```

Wall time: 4.26 s

```
Out[118]: GridSearchCV(cv=5, estimator=DecisionTreeClassifier(random_state=0),
    param_grid={'max_depth': [4, 6, 8, None],
    'min_samples_leaf': [2, 5, 1],
    'min_samples_split': [2, 4, 6]},
    refit='roc_auc',
    scoring={'accuracy', 'roc_auc', 'precision', 'f1', 'recall'})
```

```
In [119]: 1 ## Now Lets identify the optimal values for the decision tree parameter
2 tree1.best_params_
```

```
Out[119]: {'max_depth': None, 'min_samples_leaf': 5, 'min_samples_split': 2}
```

```
In [120]: 1 ##Let's identify the best AUC score achieved by the decision tree model
2 tree1.best_score_
```

```
Out[120]: 0.9817669232989432
```

This score shows that the model can predict the employees who will leave very well.

Now lets write a function that will extract all the scores from the grid search.

```
In [121]: 1 def make_results(model_name:str, model_object, metric:str):
2         '''
3         Arguments:
4             model_name (string): The model to be called in the output table
5             model_object: a fit GridSearchCV object
6             metric (string): precision, recall, f1, accuracy, or auc
7
8         Returns a pandas df with the F1, recall, precision, accuracy, and a
9         for the model with the best mean 'metric' score across all validati
10        '''
11
12        # Let's create dictionary that maps input metric to actual metric n
13        metric_dict = {'auc': 'mean_test_roc_auc',
14                       'precision': 'mean_test_precision',
15                       'recall': 'mean_test_recall',
16                       'f1': 'mean_test_f1',
17                       'accuracy': 'mean_test_accuracy'
18                       }
19
20        #Let's get all the results from the CV and put them in a df
21        cv_results = pd.DataFrame(model_object.cv_results_)
22
23        # Isolate the row of the df with the max(metric) score
24        best_estimator_results = cv_results.iloc[cv_results[metric_dict[met
25
26        # Extract Accuracy, precision, recall, and f1 score from that row
27        auc = best_estimator_results.mean_test_roc_auc
28        f1 = best_estimator_results.mean_test_f1
29        recall = best_estimator_results.mean_test_recall
30        precision = best_estimator_results.mean_test_precision
31        accuracy = best_estimator_results.mean_test_accuracy
32
33        #Let's create table of results
34        table = pd.DataFrame()
35        table = pd.DataFrame({'model': [model_name],
36                              'precision': [precision],
37                              'recall': [recall],
38                              'F1': [f1],
39                              'accuracy': [accuracy],
40                              'auc': [auc]
41                              })
42
43        return table
```

```
In [123]: 1 tree1_cv_results = make_results('Decision tree cv', tree1, 'auc')
2         tree1_cv_results
```

```
Out[123]:
```

	model	precision	recall	F1	accuracy	auc
0	Decision tree cv	0.95304	0.934866	0.943722	0.973242	0.981767

All of these scores from the decision tree model are strong indicators of good model performance.

Note That decision trees can be vulnerable to overfitting. Random forest avoid overfitting by incorporating mutiple trees to make predictions, lets develop a Random forest model

```
In [131]: 1 rf = RandomForestClassifier(random_state=0)
2         ## Let's assign a dictionary of hyperparameters to search over
3
4         cv_params = {'max_depth': [2,5,None], 'max_features':[1.0], 'max_samples'
5                     'min_samples_split':[2,3,4], 'n_estimators':[300,500]}
6
7         ##Let's instantiate GridSearch
8         rf1 = GridSearchCV(rf,cv_params,scoring=scoring,cv=5,refit='roc_auc')
```

```
In [132]: 1 %%time
2         ##Let's fit the model.
3
4         rf1.fit(X_train,y_train)
```

C:\Users\engmo\anaconda3\lib\site-packages\sklearn\metrics_classification.py:1245: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 due to no predicted samples. Use `zero_division` parameter to control this behavior.

_warn_prf(average, modifier, msg_start, len(result))

C:\Users\engmo\anaconda3\lib\site-packages\sklearn\metrics_classification.py:1245: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 due to no predicted samples. Use `zero_division` parameter to control this behavior.

_warn_prf(average, modifier, msg_start, len(result))

C:\Users\engmo\anaconda3\lib\site-packages\sklearn\metrics_classification.py:1245: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 due to no predicted samples. Use `zero_division` parameter to control this behavior.

_warn_prf(average, modifier, msg_start, len(result))

C:\Users\engmo\anaconda3\lib\site-packages\sklearn\metrics_classification.py:1245: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 due to no predicted samples. Use `zero_division` parameter to control this behavior.

```
In [136]: 1  ## Lets write a pickle functions to save and load the model results when
2  import pickle
3  path = r'C:\Users\engmo\OneDrive\Desktop\Model Pickles'
4  def write_pickle(path, model_object, save_as:str):
5      '''
6      In:
7          path:          path of folder where to save the pickle
8          model_object:  a model to pickle
9          save_as:       filename for how to save the model
10
11      Out: A call to pickle the model in the folder indicated
12      '''
13
14      with open(path + save_as + '.pickle', 'wb') as to_write:
15          pickle.dump(model_object, to_write)
16
17  def read_pickle(path, saved_model_name:str):
18
19      with open(path + saved_model_name + '.pickle', 'rb') as to_read:
20          model = pickle.load(to_read)
21
22      return model
23
```

```
In [138]: 1  ##Lets save the model in Local drive.
2  write_pickle(path, rf1, 'hr_rf1')
```

```
In [139]: 1  ## Lets read the pickle into the environment
2  rf1 = read_pickle(path, 'hr_rf1')
```

```
In [140]: 1  #Lets determine the best score
2  rf1.best_score_
```

Out[140]: 0.9906091306922387

```
In [141]: 1  ## Let's identify the best params
2  rf1.best_params_
```

Out[141]: {'max_depth': None,
'max_features': 1.0,
'max_samples': 0.7,
'min_samples_leaf': 3,
'min_samples_split': 2,
'n_estimators': 300}

```
In [142]: 1  ##Lets gather all the scores:
2  rf1_cv_results = make_results('Random Forest CV', rf1, 'auc')
3  print(tree1_cv_results)
4  print(rf1_cv_results)
```

	model	precision	recall	F1	accuracy	auc
0	Decision tree cv	0.95304	0.934866	0.943722	0.973242	0.981767
	model	precision	recall	F1	accuracy	auc
0	Random Forest CV	0.988597	0.925613	0.955976	0.979554	0.990609

The evaluation scores of the random forest model are better than those of the decision tree model, with the exception of recall (the recall score of the random forest model is lower, which is a negligible amount). This indicates that the random forest model mostly

outperforms the decision tree model.

Finally, lets evaluate the model on the test data.

```
In [143]: 1 ##Let's define a functions to get all the scores from the model's predi
2 def get_scores(model_name:str, model, X_test_data, y_test_data):
3     '''
4     Generate a table of test scores.
5
6     In:
7     model_name (string): How you want your model to be named in th
8     model: A fit GridSearchCV object
9     X_test_data: numpy array of X_test data
10    y_test_data: numpy array of y_test data
11
12    Out: pandas df of precision, recall, f1, accuracy, and AUC scores f
13    '''
14
15    preds = model.best_estimator_.predict(X_test_data)
16
17    auc = roc_auc_score(y_test_data, preds)
18    accuracy = accuracy_score(y_test_data, preds)
19    precision = precision_score(y_test_data, preds)
20    recall = recall_score(y_test_data, preds)
21    f1 = f1_score(y_test_data, preds)
22
23    table = pd.DataFrame({'model': [model_name],
24                          'precision': [precision],
25                          'recall': [recall],
26                          'f1': [f1],
27                          'accuracy': [accuracy],
28                          'AUC': [auc]
29                          })
30
31    return table
```

```
In [145]: 1 #predictions on test data
2 rf1_test_scores = get_scores('random forest test', rf1, X_test, y_test)
3 rf1_test_scores
```

```
Out[145]:
```

	model	precision	recall	f1	accuracy	AUC
0	random forest test	0.9807	0.935558	0.957597	0.9808	0.965002

This appears to be a strong model. a good indictive that it will perform good on the unseen data. We could stop here but there might be a data leakage, as explained earlier, For example the column average_monthly_hours could be a source of a leakage that the employee decided on leaving and therefore worked minimum hours.

Lets Alter the features in this model and compare the results.

```
In [146]: 1 #Lets create a feature called overworked by assigning the employees who
2 data_f = data_two.drop('satisfaction_level',axis=1)
```

In [147]:

```
1 data_f.head()
```

Out[147]:

	last_evaluation	number_project	average_monthly_hours	tenure	work_accident	left	prom
0	0.53	2	157	3	0	1	
1	0.86	5	262	6	0	1	
2	0.88	7	272	4	0	1	
3	0.87	5	223	5	0	1	
4	0.52	2	159	3	0	1	

In [148]:

```
1 #Lets add the new column.
2 data_f['overworked'] =data_f['average_monthly_hours']
3 ## as stated earlier the normal working hours per month is 166.67 hrs/m
4 data_f['overworked'] = (data_f['average_monthly_hours'] >166.67).astype(int)
5 data_f['overworked'].head()
```

Out[148]:

```
0 0
1 1
2 1
3 1
4 0
Name: overworked, dtype: int32
```

In [149]:

```
1 ## Let's drop the average_monthly_hours column.
2 data_f.drop('average_monthly_hours', axis=1,inplace=True)
3 data_f.head()
```

Out[149]:

	last_evaluation	number_project	tenure	work_accident	left	promotion_last_5years	salary
0	0.53	2	3	0	1	0	0
1	0.86	5	6	0	1	0	1
2	0.88	7	4	0	1	0	1
3	0.87	5	5	0	1	0	0
4	0.52	2	3	0	1	0	0

In [150]:

```
1 #Let's Isolate the outcome variable
2 y = data_f['left']
3
4 #selecting the features
5 X= data_f.drop('left',axis=1)
6
```

In [151]:

```
1 #lets split data
2 X_train,X_test,y_train,y_test = train_test_split(X,y,test_size=0.25,stra
```

```
In [152]: 1  ## Tree based, instantiate the model
2  tree = DecisionTreeClassifier(random_state=0)
3
4  #Assign a dictionary of hyperparameters to search over
5  cv_params = {'max_depth':[4, 6, 8, None],
6              'min_samples_leaf': [2, 5, 1],
7              'min_samples_split': [2, 4, 6]
8              }
9
10 # Assign a dictionary of scoring metrics to capture
11 scoring = {'accuracy', 'precision', 'recall', 'f1', 'roc_auc'}
12
13 # Instantiate GridSearch
14 tree2 = GridSearchCV(tree, cv_params, scoring=scoring, cv=4, refit='roc
15
```

```
In [153]: 1  %%time
2  tree2.fit(X_train, y_train)
```

Wall time: 2.96 s

```
Out[153]: GridSearchCV(cv=4, estimator=DecisionTreeClassifier(random_state=0),
    param_grid={'max_depth': [4, 6, 8, None],
    'min_samples_leaf': [2, 5, 1],
    'min_samples_split': [2, 4, 6]},
    refit='roc_auc',
    scoring={'accuracy', 'roc_auc', 'precision', 'f1', 'recall'})
```

```
In [154]: 1  #Let's check best params
2  tree2.best_params_
```

```
Out[154]: {'max_depth': None, 'min_samples_leaf': 5, 'min_samples_split': 2}
```

```
In [155]: 1  #Let's check best AUC score on CV
2  tree2.best_score_
```

```
Out[155]: 0.9707127603293837
```

This model performs very well, even without satisfaction levels and detailed hours worked data.

Next, let's check the other scores.

```
In [156]: 1  # Get all CV scores
2  tree2_cv_results = make_results('decision tree2 cv', tree2, 'auc')
3  print(tree1_cv_results)
4  print(tree2_cv_results)
```

	model	precision	recall	F1	accuracy	auc
0	Decision tree cv	0.95304	0.934866	0.943722	0.973242	0.981767
	model	precision	recall	F1	accuracy	auc
0	decision tree2 cv	0.91728	0.891332	0.904003	0.95493	0.970713

Scores fell. That's to be expected given fewer features were taken into account in this round of the model. Still, the scores are very good.

```
In [157]: 1  ## Let's see the instantiate the random forest.
2  # Instantiate model
3  rf = RandomForestClassifier(random_state=0)
4
5  # Assign a dictionary of hyperparameters to search over
6  cv_params = {'max_depth': [3,5, None],
7              'max_features': [1.0],
8              'max_samples': [0.7, 1.0],
9              'min_samples_leaf': [1,2,3],
10             'min_samples_split': [2,3,4],
11             'n_estimators': [300, 500],
12             }
13
14 # Assign a dictionary of scoring metrics to capture
15 scoring = {'accuracy', 'precision', 'recall', 'f1', 'roc_auc'}
16
17 # Instantiate GridSearch
18 rf2 = GridSearchCV(rf, cv_params, scoring=scoring, cv=4, refit='roc_auc')
```

```
In [158]: 1  %%time
2  rf2.fit(X_train, y_train)
```

```
C:\Users\engmo\anaconda3\lib\site-packages\sklearn\model_selection\_validation.py:610: FitFailedWarning: Estimator fit failed. The score on this train-test partition for these parameters will be set to nan. Details:
Traceback (most recent call last):
  File "C:\Users\engmo\anaconda3\lib\site-packages\sklearn\model_selection\_validation.py", line 593, in _fit_and_score
    estimator.fit(X_train, y_train, **fit_params)
  File "C:\Users\engmo\anaconda3\lib\site-packages\sklearn\ensemble\_forest.py", line 343, in fit
    n_samples_bootstrap = _get_n_samples_bootstrap(
  File "C:\Users\engmo\anaconda3\lib\site-packages\sklearn\ensemble\_forest.py", line 110, in _get_n_samples_bootstrap
    raise ValueError(msg.format(max_samples))
ValueError: `max_samples` must be in range (0, 1) but got value 1.0

warnings.warn("Estimator fit failed. The score on this train-test"
C:\Users\engmo\anaconda3\lib\site-packages\sklearn\model_selection\_validation.py:610: FitFailedWarning: Estimator fit failed. The score on this train-test partition for these parameters will be set to nan. Details:
```

```
In [159]: 1  # Write pickle
2  write_pickle(path, rf2, 'hr_rf2')
```

```
In [160]: 1  # Read in pickle
2  rf2 = read_pickle(path, 'hr_rf2')
```

```
In [161]: 1  # Check best params
2  rf2.best_params_
```

```
Out[161]: {'max_depth': None,
'max_features': 1.0,
'max_samples': 0.7,
'min_samples_leaf': 2,
'min_samples_split': 2,
'n_estimators': 500}
```

```
In [162]: 1 # Check best AUC score on CV
          2 rf2.best_score_
```

Out[162]: 0.9800155462940825

```
In [163]: 1 # Get all CV scores
          2 rf2_cv_results = make_results('random forest2 cv', rf2, 'auc')
          3 print(tree2_cv_results)
          4 print(rf2_cv_results)
```

	model	precision	recall	F1	accuracy	auc
0	decision tree2 cv	0.91728	0.891332	0.904003	0.95493	0.970713
	model	precision	recall	F1	accuracy	auc
0	random forest2 cv	0.933895	0.907015	0.920104	0.962486	0.980016

Again, the scores dropped slightly, but the random forest performs better than the decision tree if using AUC as the deciding metric.

Score the champion model on the test set now.

```
In [164]: 1 # Get predictions on test data
          2 rf2_test_scores = get_scores('random forest2 test', rf2, X_test, y_test)
          3 rf2_test_scores
```

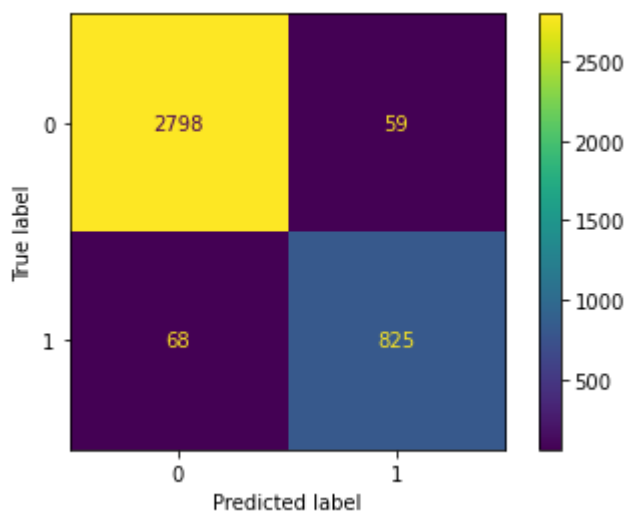
Out[164]:

	model	precision	recall	f1	accuracy	AUC
0	random forest2 test	0.933258	0.923852	0.928531	0.966133	0.951601

This seems to be a stable, well-performing final model.

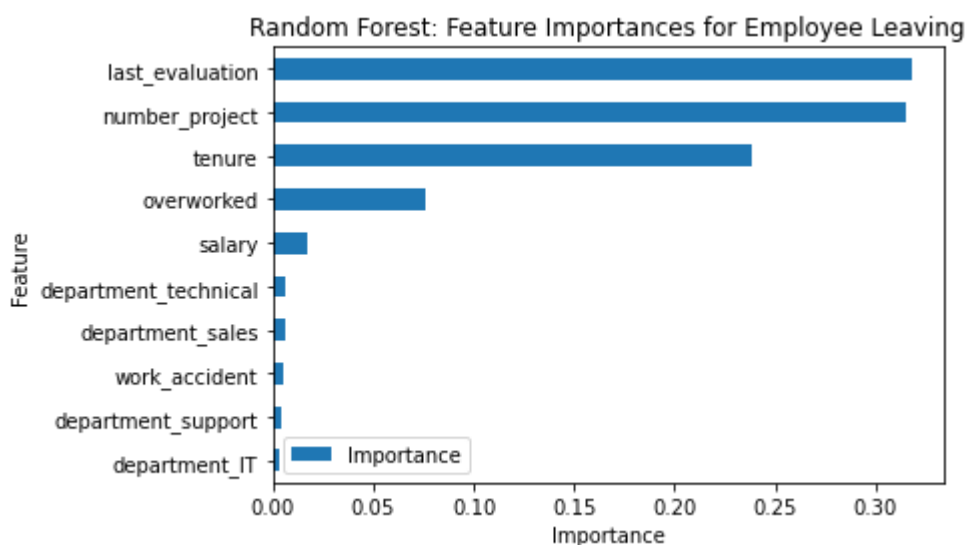
Let's plot a confusion matrix to visualize how well it predicts on the test set.

```
In [165]: 1 # Generate array of values for confusion matrix
          2 preds = rf2.best_estimator_.predict(X_test)
          3 cm = confusion_matrix(y_test, preds, labels=rf2.classes_)
          4
          5 # Plot confusion matrix
          6 disp = ConfusionMatrixDisplay(confusion_matrix=cm,
          7                             display_labels=rf2.classes_)
          8 disp.plot(values_format='');
```



For exploratory purpose, Let's inspect the most important features in the random forest model.

```
In [166]: 1 # Get feature importances
2 feat_impt = rf2.best_estimator_.feature_importances_
3
4 # Get indices of top 10 features
5 ind = np.argsort(rf2.best_estimator_.feature_importances_, -10)[-1
6
7 # Get column labels of top 10 features
8 feat = X.columns[ind]
9
10 # Filter `feat_impt` to consist of top 10 feature importances
11 feat_impt = feat_impt[ind]
12
13 y_df = pd.DataFrame({"Feature":feat,"Importance":feat_impt})
14 y_sort_df = y_df.sort_values("Importance")
15 fig = plt.figure()
16 ax1 = fig.add_subplot(111)
17
18 y_sort_df.plot(kind='barh',ax=ax1,x="Feature",y="Importance")
19
20 ax1.set_title("Random Forest: Feature Importances for Employee Leaving")
21 ax1.set_ylabel("Feature")
22 ax1.set_xlabel("Importance")
23
24 plt.show()
```



The plot above shows that in this random forest model, `last_evaluation`, `number_project`, `tenure`, and `overworked` have the highest importance, in that order. These variables are most helpful in predicting the outcome variable, `left`

Summary of model results:

- Logistic Regression:

The logistic regression model achieved precision of 81%, recall of 82%, f1-score of 81% (all weighted averages), and accuracy of 82%, on the test set.

- Tree-based Machine Learning: (Decision Tree & Random Forest)

After conducting feature engineering, the decision tree model achieved precision of 93.0%, recall of 92%, f1-score of 92%, and accuracy of 96%, on the test set. The random forest modestly outperformed the decision tree model.

Conclusion: The models and the feature importances extracted from the models confirm that employees at the company are overworked (higher hours per month & number of projects).

Recommendations:

*Cap the number of projects that employees can work on. *Consider promoting employees who have been with the company for at least four years, or conduct further investigation about why four-year tenured employees are so dissatisfied. *Either reward employees for working longer hours, or don't require them to do so. *If employees aren't familiar with the company's overtime pay policies, inform them about this. If the expectations around workload and time off aren't explicit, make them clear. *Hold company-wide and within-team discussions to understand and address the company work culture, across the board and in specific contexts. *High evaluation scores should not be reserved for employees who work 200+ hours per month. Consider a proportionate scale for rewarding employees who contribute more/put in more effort.

Next Steps:

It may be justified to still have some concern about data leakage. It could be prudent to consider how predictions change when last_evaluation is removed from the data. It's possible that evaluations aren't performed very frequently, in which case it would be useful to be able to predict employee retention without this feature. It's also possible that the evaluation score determines whether an employee leaves or stays, in which case it could be useful to pivot and try to predict performance score. The same could be said for satisfaction score.

In []:

1
